

Asymptotic Accuracy of Distribution-Based Estimation for Latent Variables

Keisuke Yamazaki

k-yam@math.dis.titech.ac.jp

Department of Computational Intelligence and Systems Science,
Tokyo Institute of Technology

G5-19 4259 Nagatsuta Midori-ku Yokohama, Japan

Abstract

Hierarchical statistical models are widely employed in information science and data engineering. The models consist of two types of variables: observable variables that represent the given data and latent variables for the unobservable labels. An asymptotic analysis of the models plays an important role in evaluating the learning process; the result of the analysis is applied not only to theoretical but also to practical situations, such as optimal model selection and active learning. There are many studies of generalization errors, which measure the prediction accuracy of the observable variables. However, the accuracy of estimating the latent variables has not yet been elucidated. For a quantitative evaluation of this, the present paper formulates distribution-based functions for the errors in the estimation of the latent variables. The asymptotic behavior is analyzed for both the maximum likelihood and the Bayes methods.

Keywords: unsupervised learning, hierarchical parametric models, latent variable, maximum likelihood method, Bayes method

1 Introduction

Hierarchical probabilistic models, such as mixture models, are mainly employed in unsupervised learning. The models have two types of variables: observable and latent. The observable variables represent the given data, and the latent ones describe the hidden data-generation process. For example, in mixture models that are employed for clustering tasks, observable variables are the attributes of the given data and the latent ones are the unobservable labels.

One of the main concerns in unsupervised learning is the analysis of the hidden processes, such as how to assign clustering labels based on the observations. Hierarchical models have an appropriate structure for this analysis, because it is straightforward to estimate the latent variables from the observable ones. Even within the limits of the clustering problem, there are a great variety of ways to detect unobservable labels, both probabilistically and deterministically, and many criteria have been proposed to evaluate the results (Dubes & Jain, 1979). For parametric models, the focus of the present paper, learning algorithms such as the expectation-maximization (EM) algorithm and the variational Bayes (VB) method (Attias, 1999; Ghahramani & Beal, 2000; Smidl & Quinn, 2005; Beal, 2003) have been developed for estimating the latent variables. These algorithms must estimate both the parameter and the variables, since the parameter is also unknown in the general case.

Theoretical analysis of the models plays an important role in evaluating the learning results. There are many studies on predicting performance in situations where both training and test data are described by the observable variables. The results of asymptotic analysis have been used for practical applications, such as model selection and active learning (Akaike, 1974; Fedorov, 1972). The simplest case of the analysis is that the learning model can attain the true model, which generates the data. Recently, it has been pointed out that when there is the redundant range/dimension of the latent variables in the learning model, singularities exist in the parameter space and the conventional statistical analysis is not valid (Amari & Ozeki, 2001). To tackle this issue, a theoretical analysis of the Bayes method was established using algebraic geometry (Watanabe, 2009). The generalization performance was then derived for various models (Yamazaki & Watanabe, 2003a; Yamazaki & Watanabe, 2003b; Rusakov & Geiger, 2005; Aoyagi, 2010; Zwiernik, 2011). Based on this analysis of the singularities, some criteria for model selection have been proposed (Watanabe, 2010; Yamazaki et al., 2005; Yamazaki et al., 2006).

Although validity of the learning algorithms is necessary for unsupervised tasks, statistical properties of the accuracy of the estimation of the latent variables have not been studied sufficiently. The goal of the present paper is to provide an asymptotic analysis for quantitative evaluation of the accuracy. For the first step, we consider the simplest case, in which the attributes, such as the range and dimension, of the latent variables are known; the true model has a minimal expression in terms of the distribution function of the observable data; and there is no singularity in the parameter space. The main contributions of the present paper are the following three items: (1) various types of estimation for the latent variables and their error functions are formulated in a distribution-based manner; (2) the asymptotic forms of the error functions are derived on the maximum likelihood and the Bayes methods; (3) it is determined that the Bayes method is more accurate than the maximum likelihood method in the asymptotic situation.

The rest of this paper is organized as follows: In Section 2 we explain the estimation of latent variables by comparing it with the prediction of observable variables. In Section 3 we provide the formal definitions of the estimation methods and the error functions. Section 4 then presents the main results for the asymptotic forms and the proofs. Discussions and conclusions are stated in Sections 5 and 6, respectively.

2 Estimations of Variables

This section distinguishes between the estimation of latent variables and the prediction of observable variables. There are variations on the estimation of latent variables due to the estimated targets.

Assume that the observable data and unobservable labels are represented by the observable variables x and the latent variables y , respectively. A set of n independent data pairs is expressed as $\{(x_1, y_1), \dots, (x_n, y_n)\}$. More precisely, there is no dependency between x_i and x_j or between y_i and y_j for $i \neq j$.

Figure 1 shows a variety of estimations of variables: prediction of an observable variable and three types of estimations of latent variables. Solid and dotted nodes are the observable and latent variables, respectively. A data pair is depicted by a connection between two nodes. The gray nodes are the target items of the estimations. We consider a stochastic approach, where the probability distribution of the target(s) is estimated from the training data $\{x_1, \dots, x_n\}$.

The top-left panel shows the prediction of unseen observable data. Based on $\{x_1, \dots, x_n\}$, the next observation $x = x_{n+1}$ is predicted. The top-right panel shows the estimation of $\{y_1, \dots, y_n\}$, which is referred to as Type I. In the stochastic approach, the joint probability of $\{y_1, \dots, y_n\}$ is estimated. The bottom-left panel shows marginal estimation, referred to as Type II. The marginal probability of y_i (y_1 is the example in the figure) is estimated; the rest of the latent variables in the probability are marginalized out. Note that there is no unseen/future data in either of Types I or II. The bottom-right panel shows estimation of y in the unseen data, which is referred to as Type III. The difference between this and Type II is the training data; the corresponding observable part of the target is included in the training set in Type II, but it is not included in Type III. In the present paper we use a distribution-based approach to analyze the theoretical accuracy of a Type-I estimation, but we also consider connections to the other types.

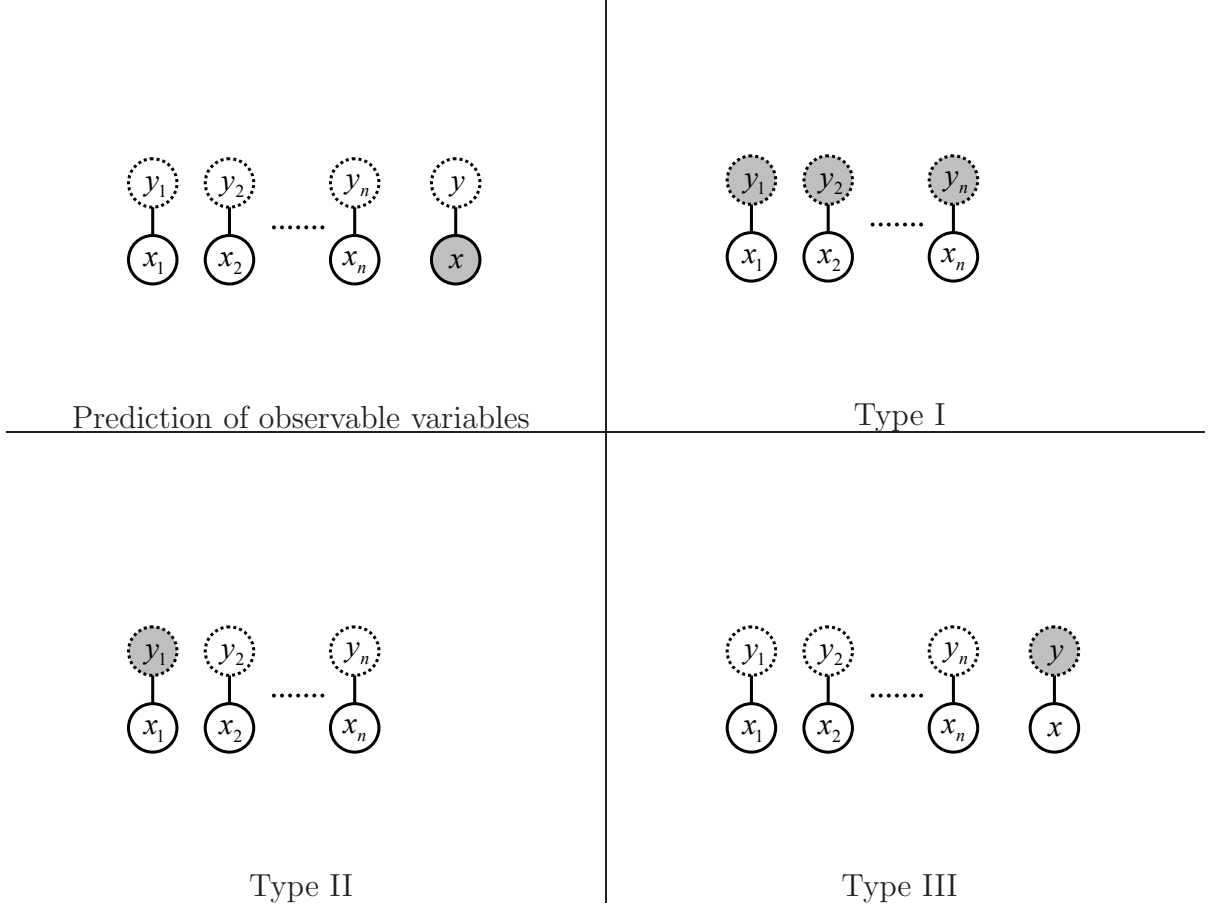


Figure 1: Prediction of observable variables and estimations of latent variables. The observable data are $\{x_1, \dots, x_n\}$. Solid and dotted nodes are observable and unobservable, respectively. Gray nodes are estimation targets.

3 Formal Definitions of Estimation Methods and Accuracy Evaluations

This section presents the maximum likelihood and Bayes methods for estimating latent variables and the corresponding error functions. Here, we consider only the Type-I estimation problem for the joint probability of the hidden part. The other types will be defined and discussed in Section 5.

Let $q(x, y) = q(y)q(x|y)$ be a joint probability of observable variables $x \in R^M$ and latent variables $y \in \{1, 2, \dots, K^*\}$. This definition indicates that both x and y are random variables and that a causal relation exists between them. In the case of a discrete x such that $x \in \{1, 2, \dots, M\}$, all the results in this paper hold if $\int dx$ is

replaced with $\sum_{x=1}^M$. The probability of the observable data x is expressed as

$$q(x) = \sum_{y=1}^{K^*} q(y)q(x|y).$$

We will refer to $q(x, y)$ as the true model.

We assume that the true probabilistic distribution with respect to x satisfies the minimality condition: the range of values of the latent variables denoted by K^* is the minimum that is required to express $q(x)$. For example, consider a three-component model, where $q(x|y=1) \neq q(x|y=2) = q(x|y=3)$ for all $x \in R^M$. The minimality condition requires the two-component expression

$$\begin{aligned} q(x) &= q(y=1)q(x|y=1) + \{q(y=2) + q(y=3)\}q(x|y=2) \\ &= q(y=1)q(x|y=1) + q(y=\bar{2})q(x|y=\bar{2}) \end{aligned}$$

where $\bar{2} = \{2, 3\}$. The minimum K^* creates a one-to-one relationship between the model expression and the function $q(x)$. If x is binary, such as $x = \{0, 1\}$ in the example, The two-component expression is redundant for describing the probabilities $q(x=0)$ and $q(x=1) = 1 - q(x=0)$. The model can be simplified to a one-component expression:

$$q(x) = q(x|y=\bar{1}),$$

where $\bar{1} = \{1, \bar{2}\}$. We find that there is no need to define the latent variables in this case, which means that the minimality condition reduces the redundancy of the latent variable space.

The notation for the data sets is $(X^n, Y^n) := \{(x_1, y_1), \dots, (x_n, y_n)\}$, $X^n = \{x_1, \dots, x_n\}$ and $Y^n = \{y_1, \dots, y_n\}$. The joint probability distribution of (X^n, Y^n) is denoted by $q(X^n, Y^n) = \prod_{i=1}^n q(x_i, y_i)$.

Let $p(x, y|w) = p(y|w)p(x|y, w)$ be a learning model, where w is the parameter and its dimension is d . Because the latent variable is unobservable, the learning model generally has its own range of variables. Then, the probability of the observable data is expressed as

$$p(x|w) = \sum_{y=1}^K p(y|w)p(x|y, w).$$

Assume that the learning model can attain the true model, i.e., there exists a set of parameters W_t such that

$$W_t = \{w^* | p(x, y|w^*) = q(x, y)\}.$$

The present paper focuses on the case $K = K^*$, where W_t consists of the unique point w^* , referred to as the true parameter. Note that there are no symmetric true parameters in W_t , because the definition is based on the joint probability distribution with respect to x and y .

We introduce two ways to construct a probability distribution of Y^n based on the observable X^n . First, we define an estimation method based on the maximum likelihood estimator. The likelihood is defined by

$$L_X(w) = \prod_{i=1}^n p(x_i|w).$$

The maximum likelihood estimator \hat{w}_X is given by

$$\hat{w}_X = \arg \max L_X(w).$$

Then, the estimated probability distribution of the latent variables is defined by

$$\begin{aligned} p(Y^n|X^n) &= \frac{p(X^n, Y^n|\hat{w}_X)}{\sum_{Y^n} p(X^n, Y^n|\hat{w}_X)} \\ &= \prod_{i=1}^n \frac{p(x_i, y_i|\hat{w}_X)}{\sum_{y_i} p(x_i, y_i|\hat{w}_X)} = \prod_{i=1}^n p(y_i|x_i, \hat{w}_X). \end{aligned} \quad (1)$$

The notation $p(Y^n|X^n, \hat{w}_X)$ is used when the method is emphasized.

Next, we define the Bayesian estimation. Let the likelihood of the joint probability distribution be

$$L_{XY}(w) = \prod_{i=1}^n p(x_i, y_i|w).$$

The marginal likelihood functions are given by

$$\begin{aligned} Z(X^n, Y^n) &= \int L_{XY}(w) \varphi(w; \eta) dw, \\ Z(X^n) &= \sum_{Y^n} Z(X^n, Y^n) = \int L_X(w) \varphi(w; \eta) dw, \end{aligned}$$

where $\varphi(w; \eta)$ is a prior with the hyperparameter η . Then, the probability of Y^n is expressed as

$$p(Y^n|X^n) = \frac{Z(X^n, Y^n)}{Z(X^n)}. \quad (2)$$

The distribution of Y^n in the true model is uniquely expressed as

$$q(Y^n|X^n) = \prod_{i=1}^n q(y_i|x_i) = \prod_{i=1}^n \frac{q(x_i, y_i)}{q(x_i)},$$

where $q(x_i) = \sum_{y_i=1}^K q(x_i, y_i)$. Accuracy of the latent variable estimation is measured by the difference between the true distribution $q(Y^n|X^n)$ and the estimated one $p(Y^n|X^n)$. For the present paper, we define the error function as the average Kullback-Leibler divergence,

$$D(n) = \frac{1}{n} E_{X^n} \left[\sum_{Y^n} q(Y^n|X^n) \ln \frac{q(Y^n|X^n)}{p(Y^n|X^n)} \right], \quad (3)$$

where the expectation is

$$E_{X^n}[f(X^n)] = \int f(X^n) q(X^n) dX^n.$$

Note that this function is available for any construction of $p(Y^n|X^n)$ when we consider the cases of the maximum likelihood and the Bayes methods below.

4 Asymptotic Analysis of the Error Function

In this section we present and prove the main theorems for the asymptotic forms of the error function.

4.1 Asymptotic Errors of the Two Methods

Let us define the following Fisher information matrices:

$$\begin{aligned} \{I_{XY}(w)\}_{ij} &= E \left[\frac{\partial \ln p(x, y|w)}{\partial w_i} \frac{\partial \ln p(x, y|w)}{\partial w_j} \right], \\ \{I_X(w)\}_{ij} &= E \left[\frac{\partial \ln p(x|w)}{\partial w_i} \frac{\partial \ln p(x|w)}{\partial w_j} \right], \end{aligned}$$

where the expectation is

$$E[f(x, y)] = \int \sum_{y=1}^K f(x, y) p(x, y|w) dx.$$

Theorem 1 *In the latent variable estimation given by Eq.1, the error function Eq.3 has the following asymptotic form:*

$$D(n) = \frac{1}{2n} \text{Tr}[\{I_{XY}(w^*) - I_X(w^*)\}I_X^{-1}(w^*)] + o\left(\frac{1}{n}\right).$$

Theorem 2 *In the latent variable estimation given by Eq.2, the error function Eq.3 has the following asymptotic form:*

$$D(n) = \frac{1}{2n} \ln \det [I_{XY}(w^*)I_X^{-1}(w^*)] + o\left(\frac{1}{n}\right).$$

These theorems reveal the speed of the decrease of the error function when the training data size n becomes large. The dominant order is $1/n$ in both methods, and its coefficient depends on the Fisher information matrices. We will present a more detailed discussion on the coefficient in Section 5.

The following corollary shows the advantage of the Bayes estimation.

Corollary 3 *Let the error functions for the maximum likelihood and the Bayes methods be denoted by $D^{ML}(n)$ and $D^{Bayes}(n)$, respectively. Assume that $I_{XY}(w^*) \neq I_X(w^*)$. For any true parameter w^* , there exists a positive constant c such that*

$$D^{ML}(n) - D^{Bayes}(n) \geq \frac{c}{n} + o\left(\frac{1}{n}\right).$$

This result shows that $D^{ML}(n) > D^{Bayes}(n)$ for a sufficiently large data size n .

4.2 Proof of Theorem 1

First, let us define another Fisher information matrix:

$$\{I_{Y|X}(w)\}_{ij} = E\left[\frac{\partial \ln p(y|x, w)}{\partial w_i} \frac{\partial \ln p(y|x, w)}{\partial w_j}\right].$$

Based on $p(y|x, w) = p(x, y|w)/p(x|w)$,

$$I_{Y|X}(w) = I_{XY}(w) + I_X(w) - J_{XY}(w) - J_{XY}^\top(w),$$

where

$$\{J_{XY}(w)\}_{ij} = E\left[\frac{\partial \ln p(x, y|w)}{\partial w_i} \frac{\partial \ln p(x|w)}{\partial w_j}\right].$$

According to the definition, we obtain

$$\begin{aligned}
\{J_{XY}(w)\}_{ij} &= E \left[\frac{1}{p(x, y|w)} \frac{\partial p(x, y|w)}{\partial w_i} \frac{\partial \ln p(x|w)}{\partial w_j} \right] \\
&= \int \sum_y \frac{\partial p(x, y|w)}{\partial w_i} \frac{\partial \ln p(x|w)}{\partial w_j} dx \\
&= \int \frac{\partial p(x|w)}{\partial w_i} \frac{\partial \ln p(x|w)}{\partial w_j} dx \\
&= \int \frac{\partial \ln p(x|w)}{\partial w_i} \frac{\partial \ln p(x|w)}{\partial w_j} p(x|w) dx = \{I_X(w)\}_{ij}.
\end{aligned}$$

Thus, it holds that

$$I_{Y|X}(w) = I_{XY}(w) - I_X(w). \quad (4)$$

Next, let us divide the error function into three parts:

$$\begin{aligned}
D(n) &= D_1(n) - D_2(n) - D_3(n), \\
D_1(n) &= \frac{1}{n} E_{X^n Y^n} [\ln q(X^n, Y^n)], \\
D_2(n) &= \frac{1}{n} E_{X^n Y^n} [\ln p(X^n, Y^n | \hat{w}_X)], \\
D_3(n) &= \frac{1}{n} E_{X^n} \left[\ln \frac{q(X^n)}{p(X^n | \hat{w}_X)} \right],
\end{aligned} \quad (5)$$

where the expectation is

$$E_{X^n Y^n} [f(X^n, Y^n)] = \int \sum_{Y^n} f(X^n, Y^n) q(X^n, Y^n) dX^n.$$

Because $D_3(n)$ is the training error on $p(x|\hat{w}_X)$, the asymptotic form is known (Akaike, 1974):

$$D_3(n) = -\frac{d}{2n} + o\left(\frac{1}{n}\right).$$

Let another estimator be defined by

$$\hat{w}_{XY} = \arg \max L_{XY}(w).$$

According to the Taylor expansion, $D_2(n)$ can be rewritten as

$$\begin{aligned}
D_2(n) &= \frac{1}{n} E_{X^n Y^n} \left[\sum_{i=1}^n \ln p(X_i, Y_i | \hat{w}_{XY}) \right] \\
&\quad + \frac{1}{n} E_{X^n Y^n} \left[\delta w^\top \sum_{i=1}^n \frac{\partial \ln p(X_i, Y_i | \hat{w}_{XY})}{\partial w} \right] \\
&\quad + \frac{1}{2n} E_{X^n Y^n} \left[\delta w^\top \sum_{i=1}^n \frac{\partial^2 \ln p(X_i, Y_i | \hat{w}_{XY})}{\partial w^2} \delta w \right] + \frac{1}{n} R_1(\delta w) \\
&= \frac{1}{n} E_{X^n Y^n} \left[\sum_{i=1}^n \ln p(X_i, Y_i | \hat{w}_{XY}) \right] \\
&\quad - \frac{1}{2} E_{X^n Y^n} [\delta w^\top I_{XY}(w^*) \delta w] + o\left(\frac{1}{n}\right),
\end{aligned}$$

where $\delta w = \hat{w}_X - \hat{w}_{XY}$, and $R_1(\delta w)$ is the remainder term. The matrix $\sum_{i=1}^n \frac{\partial^2 \ln p(X_i, Y_i | \hat{w}_{XY})}{\partial w^2}$ was replaced with $I_{XY}(w^*)$ on the basis of the law of large numbers. As for the first term of D_2 ,

$$\begin{aligned}
D_1(n) &- \frac{1}{n} E_{X^n Y^n} \left[\sum_{i=1}^n \ln p(X_i, Y_i | \hat{w}_{XY}) \right] \\
&= -\frac{d}{2n} + o\left(\frac{1}{n}\right)
\end{aligned}$$

because it is the training error on $p(x, y | \hat{w}_{XY})$. The factor in the second term of D_2 can be rewritten as

$$\begin{aligned}
&E_{X^n Y^n} [\delta w^\top I_{XY}(w^*) \delta w] \\
&= E_{X^n Y^n} [(\hat{w}_X - w^*)^\top I_{XY}(w^*) (\hat{w}_X - w^*)] \\
&\quad - E_{X^n Y^n} [(\hat{w}_{XY} - w^*)^\top I_{XY}(w^*) (\hat{w}_X - w^*)] \\
&\quad - E_{X^n Y^n} [(\hat{w}_X - w^*)^\top I_{XY}(w^*) (\hat{w}_{XY} - w^*)] \\
&\quad + E_{X^n Y^n} [(\hat{w}_{XY} - w^*)^\top I_{XY}(w^*) (\hat{w}_{XY} - w^*)].
\end{aligned} \tag{6}$$

Let us define an extended likelihood function,

$$L_2(w_{12}) = \sum_{i=1}^n \ln p(X_i, Y_i | w_1) + \sum_{i=1}^n \ln p(X_i | w_2),$$

where $w_{12} = (w_1^\top, w_2^\top)^\top$, $\hat{w}_{12} = (\hat{w}_{XY}^\top, \hat{w}_X^\top)^\top$, and $w^{**} = (w^{*\top}, w^{*\top})^\top$ are extended

vectors. According to the Taylor expansion,

$$\begin{aligned}\frac{\partial L_2(w_{12})}{\partial w_{12}} &= \left(\frac{\partial \sum \ln p(X_i, Y_i | w^*)}{\partial w_1}^\top, \frac{\partial \sum \ln p(X_i | w^*)}{\partial w_2}^\top \right)^\top \\ &\quad - M \delta w_{12}, \\ \delta w_{12} &= w_{12} - w^{**} \\ M &= \begin{bmatrix} -\frac{\partial^2 \sum \ln p(X_i, Y_i | w^*)}{\partial w_1^2} & 0 \\ 0 & -\frac{\partial^2 \sum \ln p(X_i | w^*)}{\partial w_2^2} \end{bmatrix}.\end{aligned}$$

According to $\frac{\partial L_2(\hat{w}_{12})}{\partial w_{12}} = 0$, $\delta \hat{w}_{12} = \hat{w}_{12} - w^{**}$ can be written as

$$\delta \hat{w}_{12} = M^{-1} \left(\frac{\partial \sum \ln p(X_i, Y_i | w^*)}{\partial w_1}^\top, \frac{\partial \sum \ln p(X_i | w^*)}{\partial w_2}^\top \right)^\top.$$

Based on the central limit theorem, $\delta \hat{w}_{12}$ is distributed from $\mathcal{N}(0, nM^{-1}\Sigma^{-1}M^{-1})$, where

$$\Sigma^{-1} = \begin{bmatrix} I_{XY}(w^*) & J_{XY}(w^*) \\ J_{XY}^\top(w^*) & I_X(w^*) \end{bmatrix}.$$

The covariance $nM^{-1}\Sigma^{-1}M^{-1}$ of $\delta \hat{w}_{12}$ directly shows the covariance of the estimators \hat{w}_X and \hat{w}_{XY} in Eq.6. Thus it holds that

$$\begin{aligned}& E_{X^n Y^n} [\delta w^\top I_{XY}(w^*) \delta w] \\ &= \frac{1}{n} \text{Tr} \left[I_{XY}(w^*) I_X^{-1}(w^*) \right] - \frac{1}{n} \text{Tr} \left[J_{XY}(w^*) I_X^{-1}(w^*) \right] \\ &\quad - \frac{1}{n} \text{Tr} \left[J_{XY}^\top(w^*) I_X^{-1}(w^*) \right] + \frac{1}{n} \text{Tr} \left[I_X(w^*) I_X^{-1}(w^*) \right] + o\left(\frac{1}{n}\right).\end{aligned}$$

Considering the relation Eq.5, we obtain that

$$D(n) = \frac{1}{2n} \text{Tr} [I_{Y|X}(w^*) I_X^{-1}(w^*)] + o\left(\frac{1}{n}\right).$$

Based on Eq.4, the theorem is proved. **(End of Proof)**

4.3 Proof of Theorem 2

Let us define the following entropy functions:

$$\begin{aligned}S_{XY} &= - \sum_{y=1}^{K^*} \int q(x, y) \ln q(x, y) dx, \\ S_X &= - \int q(x) \ln q(x) dx.\end{aligned}$$

According to the definition, the error function Eq.3 with the Bayes estimation can be rewritten as

$$D(n) = \frac{1}{n} \left\{ F_{XY}(n) - F_X(n) \right\},$$

where

$$\begin{aligned} F_{XY}(n) &= -nS_{XY} - E_{X^n Y^n} \left[\ln Z(X^n, Y^n) \right], \\ F_X(n) &= -nS_X - E_{X^n} \left[\ln Z(X^n) \right]. \end{aligned}$$

Based on the Taylor expansion at $w = \hat{w}_X$,

$$\begin{aligned} F_X(n) &= -nS_X - E_{X^n} \left[\ln \int \exp \left\{ \ln p(X^n | \hat{w}_X) \right. \right. \\ &\quad \left. \left. + \frac{1}{2} (w - \hat{w}_X)^\top \frac{\partial^2 \ln p(X^n | \hat{w}_X)}{\partial w^2} (w - \hat{w}_X) + r_1(w) \right\} \varphi(w; \eta) dw \right] \\ &= -nS_X - E_{X^n} [\ln p(X^n | \hat{w}_X)] - E_{X^n} \left[\ln \int e^{r_1(w)} \varphi(w; \eta) \mathcal{N}(\hat{w}_X, \Sigma_1/n) dw \right], \end{aligned}$$

where $r_1(w)$ is the remainder term and

$$\Sigma_1^{-1} = -\frac{1}{n} \frac{\partial^2 \ln p(X^n | \hat{w}_X)}{\partial w^2},$$

which converges to $I_X(w^*)$ based on the law of large numbers. Again, applying the expansion at $w = w^*$ to $e^{r_1(w)} \varphi(w; \eta)$, we obtain

$$\begin{aligned} F_X(n) &= E_{X^n} \left[\ln \frac{q(X^n)}{p(X^n | \hat{w}_X)} \right] - \ln \sqrt{2\pi^d} \sqrt{\det\{nI_X(w^*)\}^{-1}} \\ &\quad - E_{X^n} \left[\ln \int \left\{ e^{r_1(w^*)} \varphi(w^*; \eta) \right. \right. \\ &\quad \left. \left. + (w - w^*)^\top \frac{\partial e^{r_1(w^*)} \varphi(w^*; \eta)}{\partial w} + r_2(w) \right\} \mathcal{N}(\hat{w}_X, \{nI_X(w^*)\}^{-1}) dw \right] + o(1), \end{aligned}$$

where $r_2(w)$ is the remainder term. The first term is the training error on $p(x | \hat{w}_X)$. According to (Akaike, 1974), it holds that

$$E_{X^n} \left[\ln \frac{q(X^n)}{p(X^n | \hat{w}_X)} \right] = -\frac{d}{2} + o(1).$$

Then, we obtain

$$F_X(n) = \frac{d}{2} \ln \frac{n}{2\pi e} + \ln \frac{\sqrt{\det I_X(w^*)}}{\varphi(w^*; \eta)} + o(1),$$

which is consistent with the result of (Clarke & Barron, 1990). By replacing X^n with (X^n, Y^n) ,

$$F_{XY}(n) = \frac{d}{2} \ln \frac{n}{2\pi e} + \ln \frac{\sqrt{\det I_{XY}(w^*)}}{\varphi(w^*; \eta)} + o(1).$$

Therefore,

$$D(n) = \frac{1}{2n} \left\{ \ln \det I_{XY}(w^*) - \ln \det I_X(w^*) \right\} + o\left(\frac{1}{n}\right),$$

which proves the theorem. **(End of Proof)**

4.4 Proof of Corollary 3

Because $I_{XY}(w)$ is symmetric positive definite, we have a decomposition $I_{XY}(w) = LL^\top$, where L is a lower triangular matrix. The other Fisher information matrix $I_X(w)$ is also symmetric positive definite. Thus, $L^\top I_X^{-1}(w)L$ is positive definite. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$ be the eigenvalues of $L^\top I_X^{-1}(w)L$. According to the assumption, at least one eigenvalue is different from the others. Then, we obtain

$$\begin{aligned} 2n\{D^{\text{ML}}(n) - D^{\text{Bayes}}(n)\} &= \text{Tr}[I_{XY}(w)I_X^{-1}(w)] - d - \ln \det[I_{XY}(w)I_X^{-1}(w)] + o(1) \\ &= \text{Tr}[L^\top I_X^{-1}(w)L] - d - \ln \det[L^\top I_X^{-1}(w)L] + o(1) \\ &= \sum_{i=1}^d \{\lambda_i - 1\} - \ln \prod_{i=1}^d \lambda_i + o(1) \\ &= \sum_{i=1}^d \{\lambda_i - 1 - \ln \lambda_i\} + o(1). \end{aligned}$$

The first term in the last expression is positive, which proves the corollary. **(End of Proof)**

5 Discussion

5.1 Symmetry of the Learning Results

In hierarchical models, there are symmetries in both the parameter space and the latent variables. We consider the following simple case to observe them.

Example 4 Let $q(x)$ and $p(x|w)$ be Gaussian mixtures that have two components,

$$\begin{aligned} q(x) &= a^* \mathcal{N}(x; 0, 1^2) + (1 - a^*) \mathcal{N}(x; b^*, 1^2), \\ p(x|w) &= a \mathcal{N}(x; b_1, 1^2) + (1 - a) \mathcal{N}(x; b_2, 1^2), \end{aligned}$$

where $\mathcal{N}(x; \mu, \sigma^2)$ is a one-dimensional Gaussian distribution, and a^* and a are mixing ratios. The minimality condition requires that $b^* \neq 0$. The parameter of the learning model is $w = (a, b_1, b_2)$. We assume that

$$\begin{aligned} q(x, y = 1) &= a^* \mathcal{N}(x; 0, 1^2), \\ q(x, y = 2) &= (1 - a^*) \mathcal{N}(x; b^*, 1^2), \\ p(x, y = 1|w) &= a \mathcal{N}(x; b_1, 1^2), \\ p(x, y = 2|w) &= (1 - a) \mathcal{N}(x; b_2, 1^2). \end{aligned}$$

The learning model has two parameter points at which to express the true model: $w_{t1} = (a^*, 0, b^*)$ and $w_{t2} = (1 - a^*, b^*, 0)$. This is known as label switching. According to the definition of W_t , the true parameter is the unique point $w^* = (a^*, 0, b^*)$. This implies that the former expression is accepted as the proper estimation of y and the latter one, which exchanges the components, is not. This is reflected in the definition of the error function (Eq.3). The order of the components strictly eliminates the symmetries. We refer to this restriction in the error as the *asymmetric constraint*.

Let us investigate the relation between the asymmetric constraint and the error value for both methods. For a sufficiently large amount of training data, the likelihood function in Example 4 has two peaks that share the same value and are in the neighborhood of the points w_{t1} and w_{t2} . This is because there is no information on the component label from the observable data X^n . Convergence of the maximum likelihood estimator \hat{w}_X thus depends on the initial point. Theorem 1 shows the asymptotic error for the convergence to $w^* = w_{t1}$. Due to improper labeling, the estimation of w_{t2} will have a bias term in the asymptotic error, i.e., the error does not converge to zero. Therefore, Theorem 1 indicates the best performance that can be obtained by the maximum likelihood estimator, but at the same time, it indicates that the method will not always achieve this.

The Bayes estimation also has symmetries both in the latent variables and the parameter spaces. Due to the symmetry of the parameter space, the estimated distribution in Example 4 satisfies $p(Y^n|X^n) = p(\bar{Y}^n|X^n)$, where \bar{Y}^n means that labels 1 and 2 in Y^n are swapped for each other. The symmetry may adversely affect the error; the estimation result $p(Y^n|X^n)$ with parameter marginalization takes account of the symmetry, while $q(Y^n|X^n)$ does not. To more precisely investigate this effect, we eliminate the parametric symmetry and derive the asymptotic error. According to the component parameters, we divide the parameter space into two regions, such that $W_1 = (a, b_1, b_2)$ for $b_1 \leq b_2$ and W_2 for $b_1 \geq b_2$. Assume that

$b^* > 0$, where w_{t1} belongs to W_1 and w_{t2} belongs to W_2 . Let us define the distribution of Y^n as

$$p_{W_i}(Y^n|X^n) = \frac{Z_{W_i}(X^n, Y^n)}{\sum_{Y^n} Z_{W_i}(X^n, Y^n)},$$

$$Z_{W_i}(X^n, Y^n) = \int_{W_i} L_{XY}(w) \varphi(w; \eta) dw,$$

and the variants of the true model as

$$q_1(x, y) = p(x, y|w_{t1}) = q(x, y),$$

$$q_2(x, y) = p(x, y|w_{t2}).$$

We define an error function as follows:

$$D_{sym}(n) = \frac{1}{n} \min_{i=1,2} E_{X^n Y^n} \left[\ln \frac{q_i(Y^n|X^n)}{p_{W_i}(Y^n|X^n)} \right],$$

which has the true parameter in each symmetric region. We can easily obtain that $D_{sym}(n)$ is asymptotically equivalent to $D(n)$ even in the general case. Therefore, we conclude that parameter symmetry does not adversely affect the error value in the Bayes method.

5.2 Relation to Other Error Functions

We now formulate the predictions of observable data and the remaining estimations for Types II and III, and we consider the relations of their error functions to that of Type I.

First, we compare the error function to the generalization error, which measures the prediction performance on unseen observable data. The generalization error is defined as

$$D_x(n) = E_{X^n} \left[\int q(x) \ln \frac{q(x)}{p(x|X^n)} dx \right],$$

where x is independent of X^n in the data-generating process of $q(x)$. The predictive distribution $p(x|X^n)$ is constructed by

$$p(x|X^n) = p(x|\hat{w}_X)$$

for the maximum likelihood method and

$$p(x|X^n) = \int p(x|w) p(w|X^n) dw$$

for the Bayes method. The posterior distribution of the Bayes method is given by

$$p(w|X^n) = \frac{L_X(w)\varphi(w;\eta)}{Z(X^n)}.$$

Both methods have the same dominant terms in their asymptotic forms,

$$D_x(n) = \frac{d}{2n} + o\left(\frac{1}{n}\right).$$

The coefficient of the asymptotic generalization error depends only on the dimension of the parameter for any model, but that of $D(n)$ is determined by both the model expression and the true parameter w^* .

Next, we discuss Type-II estimation; we focus on the value y_i from Y^n and its estimation accuracy. Based on the joint probability, the estimation of y_i is defined by

$$p(y_i|X^n) = \sum_{Y^n \setminus y_i} p(Y^n|X^n),$$

where the summation is taken over Y^n except for y_i . Thus the error function depends on which y_i we exclude. In order to measure the average effect of the exclusions, we define the error as follows:

$$D_{y|X^n}(n) = E_{X^n} \left[\frac{1}{n} \sum_{i=1}^n \sum_{y_i} q(y_i|x_i) \ln \frac{q(y_i|x_i)}{p(y_i|X^n)} \right].$$

The maximum likelihood method has the following estimation,

$$\begin{aligned} p(y_i|X^n) &= \sum_{Y^n \setminus y_i} \prod_{i=1}^n \frac{p(x_i, y_i|\hat{w}_X)}{p(x_i|\hat{w}_X)} \\ &= \frac{p(x_1|\hat{w}_X) \cdots p(x_{i-1}|\hat{w}_X) p(x_i, y_i|\hat{w}_X) p(x_{i+1}|\hat{w}_X) \cdots p(x_n|\hat{w}_X)}{\prod_{i=1}^n p(x_i|\hat{w}_X)} \\ &= \frac{p(x_i, y_i|\hat{w}_X)}{p(x_i|\hat{w}_X)} = p(y_i|x_i, \hat{w}_X). \end{aligned}$$

We can easily find that

$$\begin{aligned} D_{y|X^n}(n) &= E_{X^n} \left[\frac{1}{n} \sum_{i=1}^n \sum_{y_i=1}^K q(y_i|x_i) \ln \frac{q(y_i|x_i)}{p(y_i|x_i, \hat{w}_X)} \right] \\ &= \frac{1}{n} E_{X^n} \left[\sum_{Y^n} q(Y^n|X^n) \ln \frac{q(Y^n|X^n)}{p(Y^n|X^n, \hat{w}_X)} \right]. \end{aligned}$$

Therefore, it holds that $D_{y|X^n}(n) = D(n)$ in the maximum likelihood method. However, the Bayes method has the estimation,

$$p(y_i|X^n) = \frac{\int p(x_1|w) \cdots p(x_{i-1}|w) p(x_i, y_i|w) p(x_{i+1}|w) \cdots p(x_n|w) \varphi(w; \eta) dw}{Z(X^n)},$$

which indicates $D_{y|X^n}(n) \neq D(n)$. A sufficient condition for $D_{y|X^n}(n) = D(n)$ is to satisfy $p(Y^n|X^n) = \prod_{i=1}^n p(y_i|X^n)$.

Finally, we consider the Type-III estimation. The error is defined by

$$D_{y|x}(n) = E_{X^n} \left[\int q(x) \sum_{y=1}^K q(y|x) \ln \frac{q(y|x)}{p(y|x, X^n)} dx \right].$$

Note that the new observation x is not used for estimation of y , or $D_{y|x}(n)$ will be equivalent to the Type-II error $D_{y|X^{n+1}}(n+1)$. The maximum likelihood estimation $p(y|x, X^n)$ is given by

$$p(y|x, X^n) = \frac{p(x, y|\hat{w}_X)}{p(x|\hat{w}_X)},$$

and for the Bayes method it is

$$p(y|x, X^n) = \int \frac{p(x, y|w)}{p(x|w)} p(w|X^n) dw. \quad (7)$$

Using the result in (Shimodaira, 1993) for a variant Akaike information criterion (AIC) from partially observed data, we immediately obtain the asymptotic form of $D_{y|x}(n)$ as

$$D_{y|x}(n) = \frac{1}{2n} \text{Tr} \left[\left\{ I_{XY}(w^*) - I_X(w^*) \right\} I_X(w^*)^{-1} \right] + o\left(\frac{1}{n}\right).$$

We thus conclude that all estimation types have the same accuracy in the maximum likelihood method. The difference of the training data between Types II and III does not asymptotically affect the estimation results. The analysis of the Type-III estimate in the Bayes method is left for future study.

5.3 Variants of Types II and III

Table 1 summarizes the results in the previous subsection. The rows indicate the maximum likelihood (ML) and the Bayes methods, respectively. The Fisher information matrices $I_{XY}(w^*)$ and $I_X(w^*)$ are abbreviated in a form that does not include the true parameter, i.e., I_{XY} and I_X . The error functions of Types II and III in the

Table 1: Coefficients of the dominant order $1/n$ in the error functions

	Prediction	Type I	Type II	Type III
ML	$d/2$	$\text{Tr}[\{I_{XY} - I_X\}I_X^{-1}]/2$	$\text{Tr}[\{I_{XY} - I_X\}I_X^{-1}]/2$	$\text{Tr}[\{I_{XY} - I_X\}I_X^{-1}]/2$
Bayes	$d/2$	$\ln \det[I_{XY}I_X^{-1}]/2$	unknown	unknown

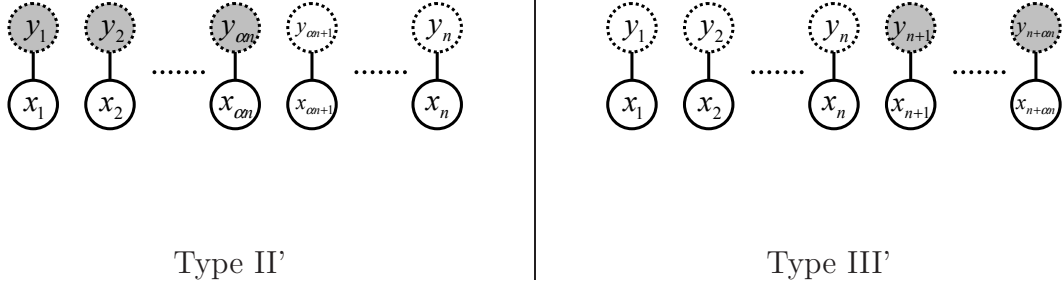


Figure 2: (Left) Partial marginal estimation for $y_1, \dots, y_{\alpha n}$. (Right) Estimation for future data $y_{n+1}, \dots, y_{n+\alpha n}$.

Bayes method are still unknown. The analysis is not straightforward when there is a single target of estimation, because the asymptotic expansion is not available when the number of target nodes is constant with respect to the training data size n .

Consider the variants of Types II and III depicted in Figure 2. Assume that $0 < \alpha \leq 1$ is a constant rational number and that n gets large enough to satisfy that αn is an integer. The left panel shows the partial marginal estimation referred to as Type II'. We will consider the joint probability of $y_1, \dots, y_{\alpha n}$, where the remaining variables $y_{\alpha n+1}, \dots, y_n$ have been marginalized out. Type II' is equivalent to Type I when $\alpha = 1$. Note that the order in which the target nodes are determined does not change the average accuracy for i.i.d. data. The right panel indicates the estimations for future data $y_{n+1}, \dots, y_{n+\alpha n}$. We refer to it as Type III' and construct the joint probability on these variables. In the variant types, the targets are changed from a single node to αn nodes, which enables us to analyze the asymptotic behavior.

We will use the following notation:

$$\begin{aligned} X_1 &= \{x_1, \dots, x_{\alpha n}\}, \\ Y_1 &= \{y_1, \dots, y_{\alpha n}\} \end{aligned}$$

for Type II' and

$$\begin{aligned} X_2 &= \{x_{n+1}, \dots, x_{n+\alpha n}\}, \\ Y_2 &= \{y_{n+1}, \dots, y_{n+\alpha n}\} \end{aligned}$$

for Type III'. The Bayes estimations are given by

$$\begin{aligned} p(Y_1|X^n) &= \frac{\int \prod_{j=1}^{\alpha n} p(x_j, y_j|w) \prod_{i=\alpha n+1}^n p(x_i|w) \varphi(w; \eta) dw}{\int \prod_{i=1}^n p(x_i|w) \varphi(w; \eta) dw}, \\ p(Y_2|X_2, X^n) &= \int \prod_{i=n+1}^{n+\alpha n} \frac{p(x_i, y_i|w)}{p(x_i|w)} p(w|X^n) dw \end{aligned}$$

for Type II' and Type III', respectively. The respective error functions are defined by

$$\begin{aligned} D_{Y_1|X^n}(n) &= \frac{1}{\alpha n} E_{X^n} \left[\sum_{Y_1} q(Y_1|X^n) \ln \frac{q(Y_1|X^n)}{p(Y_1|X^n)} \right], \\ D_{Y_2|X_2}(n) &= \frac{1}{\alpha n} E_{X^n, X_2} \left[\sum_{Y_2} q(Y_2|X_2) \ln \frac{q(Y_2|X_2)}{p(Y_2|X_2, X^n)} \right]. \end{aligned}$$

In ways similar to the proofs of Theorems 1 and 2, the asymptotic forms are derived as follows.

Theorem 5 *In Type II', the error function has the following asymptotic form:*

$$D_{Y_1|X^n}(n) = \frac{1}{2\alpha n} \ln \det[K_{XY}(w^*) I_X(w^*)^{-1}] + o\left(\frac{1}{n}\right),$$

where $K_{XY}(w) = \alpha I_{XY}(w) + (1 - \alpha) I_X(w)$.

The proof is in the appendix.

Theorem 6 *In Type III', the error function has the following asymptotic form:*

$$D_{Y_2|X_2}(n) = \frac{1}{2\alpha n} \ln \det[K_{XY}(w^*) I_X^{-1}(w^*)] + o\left(\frac{1}{n}\right).$$

This proof is also in the appendix. These theorems show that when Types II' and III' have the same α , they asymptotically have the same accuracy. This implies the asymptotic equivalency of Types II and III by combining the results of the maximum likelihood method.

Table 2 summarizes the results. Based on the definitions, the results for the maximum likelihood method are also available for Types II' and III'. Using the asymptotic forms, we can compare the relation of the magnitudes for the maximum likelihood method.

Table 2: Coefficients of the dominant order $1/n$ in the error functions

	Pred.	Type I	Type II'	Type III'
ML	$d/2$	$\text{Tr}[\{I_{XY} - I_X\}I_X^{-1}]/2$	$\text{Tr}[\{I_{XY} - I_X\}I_X^{-1}]/2$	$\text{Tr}[\{I_{XY} - I_X\}I_X^{-1}]/2$
Bayes	$d/2$	$\ln \det[I_{XY}I_X^{-1}]/2$	$\ln \det[K_{XY}I_X^{-1}]/(2\alpha)$	$\ln \det[K_{XY}I_X^{-1}]/(2\alpha)$

Corollary 7 *Assume that $I_{XY}(w) \neq I_X(w)$. For $0 < \alpha \leq 1$, there exists a positive constant c_1 such that*

$$\text{Tr}[\{I_{XY}(w) - I_X(w)\}I_X^{-1}(w)] - \frac{1}{\alpha} \ln \det[K_{XY}(w)I_X^{-1}(w)] \geq \frac{c_1}{n} + o\left(\frac{1}{n}\right).$$

The proof is in the appendix. We immediately obtain the following relation, which shows the advantage of the Bayes estimation in the asymptotic case:

$$\begin{aligned} D_{Y_1|X^n}^{\text{Bayes}}(n) &< D_{Y_1|X^n}^{\text{ML}}(n) \\ D_{Y_2|X_2}^{\text{Bayes}}(n) &< D_{Y_2|X_2}^{\text{ML}}(n) \end{aligned}$$

for respective α 's.

By comparing the errors of Types I and II' in the Bayes method, we can obtain the effect of supplementary observable data. Let us consider the Type-II' case in which the estimation target is Y_1 and the training data is only X_1 . This corresponds to the estimation in Type I with αn training data, which we emphasize by calling it Type I'. The difference between Type I' and Type II' is the addition of supplementary data $X^n \setminus X_1$.

Corollary 8 *Assume that the minimum eigenvalue of $I_{XY}(w^*)I_X^{-1}(w^*)$ is not less than one, i.e., $\lambda_d \geq 1$. The error difference is asymptotically described as*

$$\begin{aligned} D(\alpha n) - D_{Y_1|X^n}(n) &= \frac{1}{2\alpha n} \ln \det[I_{XY}(w^*)K_{XY}^{-1}(w^*)] + o\left(\frac{1}{n}\right) \\ &\geq \frac{c_2}{n} + o\left(\frac{1}{n}\right), \end{aligned}$$

where c_2 is a positive constant. This shows that Type II' has a smaller error than Type I' in the asymptotic situation; the supplementary data make the estimation more accurate.

The proof is in the appendix.

5.4 Comparison between the Two Methods

Corollaries 3 and 7 show that the Bayes method is more accurate than the maximum likelihood method for Types I, II', and III'. There have been many data-based

comparisons of the predicting performances of these two methods (e.g., (Akaike, 1980; ?; ?)). We will now discuss the computational costs of the two methods for the estimation of latent variables. We note there will be a trade-off between cost and accuracy.

We will assume that the estimated distribution is to be calculated for a practical purpose. For example, the value of $p(Y^n|X^n)$ in Type I is used for sampling label assignments and for searching for the optimal assignment $\arg \max_{Y^n} p(Y^n|X^n)$. The maximum likelihood method requires the determination of \hat{w}_X for all Types I, II, and III. The computation is not expensive once \hat{w}_X is successfully found, but the global maximum point of the likelihood function is not easily obtained. The EM algorithm is commonly used for searching for the maximum likelihood estimator in models with latent variables, but it is often trapped in one of the local maxima. The results of the steepest descent method also depend on the initial point and the step size of the iteration.

The Bayes method is generally expensive. In the estimated distribution $p(Y^n|X^n)$ of Type I, the numerator $Z(X^n, Y^n)$ contains integrals that depend on Y^n . Sampling y_i in Type II requires the same computation as for Type I: we can obtain y_i by ignoring the other elements $Y^n \setminus y_i$, which realizes the marginalization $\sum_{Y^n \setminus y_i} p(Y^n|X^n)$. A conjugate prior allows us to have a tractable form of $Z(X^n, Y^n)$ (Dawid & Lauritzen, 1993; Heckerman, 1999), which reduces the computational cost. In Type III, Eq.7 shows that there is no direct sampling method for y . In this case, expensive sampling from the posterior $p(w|X^n)$ is necessary.

The VB method is an approximation that allows the direct computation of $P(Y^n|X^n)$ and $p(w|X^n)$, which have tractable forms and reduced computational costs. However, the assumption that $P(Y^n|X^n)$ and $p(w|X^n)$ are independent does not hold in many cases. We conjecture that the $P(Y^n|X^n)$ of the VB method will be less accurate than that of the original Bayes method.

6 Conclusions

In the present paper we formalized the estimation from the observable data of the distribution of the latent variables, and we measured its accuracy by using the Kullback-Leibler divergence. We succeeded in deriving the asymptotic error functions for both the maximum likelihood and the Bayes methods. These results allow us to mathematically compare the estimation methods: we determined that the Bayes method is more accurate than the maximum likelihood method in most cases, while their prediction accuracies are equivalent. The generalization error has been approximated from the given observable data, such as by using the cross-validation and bootstrap methods, but there is no approximation technique for the error of the estimation of the latent variables, because the latent data can not be obtained.

Therefore, these asymptotic forms are thus far the only way we have to estimate their accuracy.

Acknowledgement

This research was partially supported by the Kayamori Foundation of Informational Science Advancement and KAKENHI 23500172.

Appendix

In this section, we prove Theorem 5, Theorem 6, Corollary 7, and Corollary 8.

Proof of Theorem 5

The error function is rewritten as

$$\begin{aligned} D_{Y_1|X^n}(n) &= \frac{1}{\alpha n} \left\{ F_{XY}^{(1)}(n) - F_X(n) \right\}, \\ F_{XY}^{(1)}(n) &= -\alpha n S_{XY} - (1 - \alpha)n S_X - E_{X^n, Y_1} \left[\ln \int L_{XY}^{(1)}(w) \varphi(w; \eta) dw \right], \\ L_{XY}^{(1)}(w) &= \prod_{j=1}^{\alpha n} p(x_j, y_j | w) \prod_{i=\alpha n+1}^n p(x_i | w). \end{aligned}$$

Based on the Taylor expansion at $w = \hat{w}^{(1)}$, where $\hat{w}^{(1)} = \arg \max L^{(1)}(w)$,

$$\begin{aligned} F_{XY}^{(1)}(n) &= E_{X^n, Y_1} \left[\sum_{j=1}^{\alpha n} \ln \frac{q(x_j, y_j)}{p(x_j, y_j | \hat{w}^{(1)})} + \sum_{i=\alpha n+1}^n \ln \frac{q(x_i)}{p(x_i | \hat{w}^{(1)})} \right. \\ &\quad \left. + \ln \int \exp \left\{ -n(w - \hat{w}^{(1)})^\top G^{(1)}(X^n, Y_1)(w - \hat{w}^{(1)}) + r_3(w) \right\} \varphi(w; \eta) dw \right], \end{aligned}$$

where $r_3(w)$ is the remainder term and

$$G^{(1)}(X^n, Y_1) = -\frac{1}{n} \frac{\partial^2}{\partial w^2} \left(\sum_{j=1}^{\alpha n} \ln p(x_j, y_j | \hat{w}^{(1)}) + \sum_{i=\alpha n+1}^n \ln p(x_i | \hat{w}^{(1)}) \right).$$

The first and the second terms of $F_{XY}^{(1)}(n)$ correspond to the training error. Following the same method as we used in the proof of Theorem 2 and noting that

$$G^{(1)}(X^n, Y_1) \rightarrow K_{XY}(w^*),$$

we obtain

$$F_{XY}^{(1)}(n) = \frac{d}{2} \ln \frac{n}{2\pi e} + \ln \frac{\sqrt{\det K_{XY}(w^*)}}{\varphi(w^*; \eta)} + o(1),$$

which completes the proof. **(End of Proof)**

Proof of Theorem 6

The error function is rewritten as

$$\begin{aligned} D_{Y_2|X_2}(n) &= \frac{1}{\alpha n} \left\{ F_{XY}^{(2)}(n) - F_X(n) \right\}, \\ F_{XY}^{(2)}(n) &= -\alpha n S_{XY} - n S_X - E_{X^n, X_2, Y_2} \left[\ln \int L_{XY}^{(2)}(w) \varphi(w; \eta) dw \right], \\ L_{XY}^{(2)}(w) &= \prod_{j=n+1}^{n+\alpha n} p(y_j | x_j, w) \prod_{i=1}^n p(x_i | w). \end{aligned}$$

Based on the Taylor expansion at $w = \hat{w}^{(2)}$, where $\hat{w}^{(2)} = \arg \max L^{(2)}(w)$,

$$\begin{aligned} F_{XY}^{(2)}(n) &= E_{X^n, X_2, Y_2} \left[\sum_{j=n+1}^{\alpha n} \ln \frac{q(y_j | x_j)}{p(y_j | x_j, \hat{w}^{(2)})} + \sum_{i=1}^n \ln \frac{q(x_i)}{p(x_i | \hat{w}^{(2)})} \right. \\ &\quad \left. + \ln \int \exp \left\{ -n(w - \hat{w}^{(2)})^\top G^{(2)}(X^n, X_2, Y_2)(w - \hat{w}^{(2)}) + r_4(w) \right\} \varphi(w; \eta) dw \right], \end{aligned}$$

where $r_4(w)$ is the remainder term and

$$G^{(2)}(X^n, X_2, Y_2) = -\frac{1}{n} \frac{\partial^2}{\partial w^2} \left(\sum_{j=n+1}^{\alpha n} \ln p(y_j | x_j, \hat{w}^{(2)}) + \sum_{i=1}^n \ln p(x_i | \hat{w}^{(2)}) \right).$$

The first and the second terms of $F_{XY}^{(1)}(n)$ correspond to the training error, which are stated as

$$\begin{aligned} E_{X^n, X_2, Y_2} \left[\sum_{j=n+1}^{\alpha n} \ln \frac{q(y_j | x_j)}{p(y_j | x_j, \hat{w}^{(2)})} + \sum_{i=1}^n \ln \frac{q(x_i)}{p(x_i | \hat{w}^{(2)})} \right] \\ = -\text{Tr} \left[\{ \alpha I_{Y|X}(w^*) + I_X(w^*) \} K_{XY}(w^*)^{-1} \right] + o(1). \end{aligned}$$

Following the same method we used in the proof of Theorem 2 and noting that

$$G^{(2)}(X^n, X_2, Y_2) \rightarrow K_{XY}(w^*),$$

we obtain

$$\begin{aligned}
F_{XY}^{(1)}(n) &= -\text{Tr} \left[\{ \alpha I_{Y|X}(w^*) + I_X(w^*) \} K_{XY}(w^*)^{-1} \right] \\
&\quad + \frac{d}{2} \ln \frac{n}{2\pi} + \ln \frac{\sqrt{\det K_{XY}(w^*)}}{\varphi(w^*; \eta)} + o(1) \\
&= \frac{d}{2} \ln \frac{n}{2\pi e} + \ln \frac{\sqrt{\det K_{XY}(w^*)}}{\varphi(w^*; \eta)} + o(1),
\end{aligned}$$

which completes the proof. **(End of Proof)**

Proof of Corollary 7

It holds that

$$\frac{1}{\alpha} \ln \det[K_{XY}(w) I_X^{-1}(w)] = \frac{1}{\alpha} \ln \det[\alpha \{I_{XY}(w) - I_X(w)\} I_X^{-1}(w) + E_d],$$

where E_d is the $d \times d$ unit matrix. On the other hand,

$$\text{Tr}[\{I_{XY}(w) - I_X(w)\} I_X^{-1}(w)] = \frac{1}{\alpha} \left\{ \text{Tr}[\alpha \{I_{XY}(w) - I_X(w)\} I_X^{-1}(w) + E_d] - d \right\}.$$

It is easy to confirm that $\alpha L_1^\top I_X^{-1}(w) L_1 + E_d$ is positive definite, where $L_1^\top L_1 = I_{XY}(w) - I_X(w)$. Considering the eigenvalues $\mu_1 \geq \mu_2 \geq \dots \geq \mu_d > 0$, we can obtain the following relation in the same way as we did in the proof of Corollary 3:

$$\text{Tr}[\{I_{XY}(w) - I_X(w)\} I_X^{-1}(w)] - \frac{1}{\alpha} \ln \det[K_{XY}(w) I_X^{-1}(w)] = \frac{1}{\alpha} \sum_{i=1}^d \left\{ \mu_i - 1 - \ln \mu_i \right\}.$$

It is easy to confirm that the right-hand side is positive, which completes the proof. **(End of Proof)**

Proof of Corollary 8

Based on the eigenvalues of $I_{XY}(w^*) I_X^{-1}(w^*)$, it holds that

$$\begin{aligned}
\ln \det[I_{XY}(w^*) K_{XY}^{-1}(w^*)] &= \ln \det[I_{XY}(w^*) I_X^{-1}(w^*)] - \ln \det[\alpha I_{XY}(w^*) I_X^{-1}(w^*) + (1 - \alpha) E_d] \\
&= \sum_{i=1}^d \ln \lambda_i - \sum_{i=1}^d \ln \{ \alpha \lambda_i + (1 - \alpha) \} \geq 0,
\end{aligned}$$

which completes the proof. **(End of Proof)**

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Trans. on Automatic Control*, 19, 716–723.
- Akaike, H. (1980). Likelihood and bayes procedure. *J. M. Bernald, Bayesian statistics* (pp. 143–166). Valencia, Italy: University Press.
- Amari, S., & Ozeki, T. (2001). Differential and algebraic geometry of multilayer perceptrons. *IEICE Trans, E84-A 1*, 31–38.
- Aoyagi, M. (2010). Stochastic complexity and generalization error of a restricted boltzmann machine in bayesian estimation. *Journal of Machine Learning Research*, 11, 1243–1272.
- Attias, H. (1999). Inferring parameters and structure of latent variable models by variational Bayes. *Proceedings of Uncertainty in Artificial Intelligence*.
- Beal, M. J. (2003). *Variational algorithms for approximate bayesian inference* (Technical Report).
- Clarke, B., & Barron, A. R. (1990). Information-theoretic asymptotics of bayes methods. *IEEE Transactions on Information Theory*, 36, 453–471.
- Dawid, A. P., & Lauritzen, S. L. (1993). Hyper-Markov laws in the statistical analysis of decomposable graphical models. *Annals of Statistics*, 21, 1272–1317.
- Dubes, R., & Jain, A. K. (1979). Validity studies in clustering methodologies. *Pattern Recognition*, 11, 235–254.
- Fedorov, V. V. (1972). *Theory of optimal experiments*. New York: Academic Press.
- Ghahramani, Z., & Beal, M. J. (2000). Graphical models and variational methods. *Advanced Mean Field Methods - Theory and Practice*. MIT Press.
- Heckerman, D. (1999). Learning in graphical models. chapter A tutorial on learning with Bayesian networks, 301–354. Cambridge, MA, USA: MIT Press.
- Rusakov, D., & Geiger, D. (2005). Asymptotic model selection for naive bayesian networks. *Journal of Machine Learning Research*, 6, 1–35.
- Shimodaira, H. (1993). A new criterion for selecting models from partially observed data. *Oldford, Eds., Selecting Models from Data: Artificial Intelligence and Statistics IV, Lecture Notes in Statistics 89* (pp. 381–386). Springer-Verlag.

- Smidl, V., & Quinn, A. (2005). *The variational bayes method in signal processing (signals and communication technology)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc.
- Watanabe, S. (2009). *Algebraic geometry and statistical learning theory*. New York, NY, USA: Cambridge University Press.
- Watanabe, S. (2010). Equations of states in singular statistical estimation. *Neural Networks*, 23, 20–34.
- Yamazaki, K., Nagata, K., & Watanabe, S. (2005). A new method of model selection based on learning coefficient. *Proceedings of International Symposium on Nonlinear Theory and its Applications* (pp. 389–392).
- Yamazaki, K., Nagata, K., Watanabe, S., & Müller, K.-R. (2006). A model selection method based on bound of learning coefficient. *LNCs* (pp. 371–380). Springer.
- Yamazaki, K., & Watanabe, S. (2003a). Singularities in mixture models and upper bounds of stochastic complexity. *International Journal of Neural Networks*, 16, 1029–1038.
- Yamazaki, K., & Watanabe, S. (2003b). Stochastic complexity of bayesian networks. *Proc. of UAI* (pp. 592–599).
- Zwiernik, P. (2011). An asymptotic behaviour of the marginal likelihood for general markov models. *J. Mach. Learn. Res.*, 999888, 3283–3310.